



Shri Vile Parle Kelavani Mandal's

DWARKADAS J. SANGHVI COLLEGE OF ENGINEERING

(Autonomous College Affiliated to the University of Mumbai)

NAAC Accredited with "A" Grade (CGPA : 3.18)



Academic Year (2022-23)

Year: 3 Semester: V

Program: Minor in Data Science

Subject: Foundation of Data Analysis

Date: 3/01/2023

Max. Marks: 75

Time: 10: 30 am to 1:30 pm

Duration: 3 Hours

REGULAR EXAMINATION

Instructions: Candidates should read carefully the instructions printed on the question paper and on the cover page of the Answer Book, which is provided for their use.

- (1) This question paper contains 03 pages.
- (2) **All Questions are Compulsory.**
- (3) All questions carry equal marks.
- (4) **Answer to each new question is to be started on a fresh page.**
- (5) **Figures in the brackets on the right indicate full marks.**
- (6) **Assume suitable data wherever required, but justify it.**
- (7) **Draw the neat labelled diagrams, wherever necessary.**

Question No.		Max. Marks																
Q1 (a)	<p>i. Calculate the arithmetic mean for the following weight recorded to the nearest grams of 60 apples picked out at random from a consignment are given below:</p> <table><tr><th>Weights (grams)</th><th>Frequency</th></tr><tr><td>65-84</td><td>09</td></tr><tr><td>85-104</td><td>10</td></tr><tr><td>105-124</td><td>17</td></tr><tr><td>125-144</td><td>10</td></tr><tr><td>145-164</td><td>05</td></tr><tr><td>165-184</td><td>04</td></tr><tr><td>185-204</td><td>05</td></tr></table> <p style="text-align: center;">OR</p> <p>ii. Answer the following:</p> <p>a. Calculate mode for ungrouped data. $X_i = 2\ 3\ 8\ 4\ 6\ 3\ 2\ 5\ 3$</p> <p>b. Calculate the median for the following grade points obtained by 10 practitioners are given. $X_i = 45\ 32\ 37\ 46\ 39\ 36\ 41\ 48\ 36\ 50$</p>	Weights (grams)	Frequency	65-84	09	85-104	10	105-124	17	125-144	10	145-164	05	165-184	04	185-204	05	<p>[05]</p> <p>[03]</p> <p>[02]</p>
Weights (grams)	Frequency																	
65-84	09																	
85-104	10																	
105-124	17																	
125-144	10																	
145-164	05																	
165-184	04																	
185-204	05																	
Q1 (b)	<p>i. Explain the relative measures of dispersion in statistics.</p> <p>ii. Describe the steps used to calculate the score for a given percentile P.</p> <p>iii. The scores of students are 3, 5, 7, 8, 9, 11, 13, 15. What is the score of the 25th percentile?</p>	<p>[05]</p> <p>[02]</p> <p>[03]</p>																
Q2 (a)	<p>1. Answer the following:</p> <p>i. Explain in detail the importance of data cube.</p> <p>ii. Sketch a 3D data cube for sales of an organization.</p> <p>iii. Sketch a Star schema of sales where Dealer, Model, Date, Product, Branch are dimension table & revenue is a fact table with Units sold and total revenue. Assume necessary attributes.</p> <p style="text-align: center;">OR</p>	<p>[03]</p> <p>[03]</p> <p>[04]</p>																



	<div>2. Answer the following:</div> <div><div>i. Explain various OLAP models in data warehouse.</div><div>ii. Consider the following data illustrating temperature of certain weeks:</div></div> <table><tr><th>Temperature</th><th>cool</th><th>mild</th><th>hot</th></tr><tr><th>Week1</th><td>2</td><td>1</td><td>1</td></tr><tr><th>Week2</th><td>2</td><td>1</td><td>1</td></tr></table> <div>Suppose we want to set daily temperature from the above data. Specify the operation used to create a new dataset.</div>	Temperature	cool	mild	hot	Week1	2	1	1	Week2	2	1	1	<div>[05]</div> <div>[05]</div>
Temperature	cool	mild	hot											
Week1	2	1	1											
Week2	2	1	1											
Q2.(b)	<div>Select the appropriate answers for the below questions and justify your answer.</div> <div><div>i. Data visualization tools provide an accessible way to see and understand . _____ in data.</div><div><div>a) trends</div><div>b) outliers</div><div>c) patterns</div><div>d) All of the above</div></div><div>ii. The charts that are helpful in making comparisons with features are _____.</div><div><div>a) Bar charts</div><div>b) Column charts</div><div>c) Pie charts</div><div>d) Both bar and column charts</div></div><div>iii. A data visualization that updates in real-time and gives multiple outputs is called as _____.</div><div><div>a) Dashboard</div><div>b) Metrics table</div><div>c) Data table</div><div>d) None of the above</div></div><div>iv. For creating variable size bins we use _____.</div><div><div>a) Sets</div><div>b) Groups</div><div>c) Calculated Fields</div><div>d) Table Calculations</div></div><div>v. Which graphs are generally used to show the relationship among the variables?</div><div><div>a) Bar graph</div><div>b) Line graph</div><div>c) Scatter Plot</div><div>d) Maps</div></div></div>	<div>[01]</div> <div>[01]</div> <div>[01]</div> <div>[01]</div> <div>[01]</div>												
Q3 (a)	<div>i. What is the need of data preprocessing? Explain the various techniques of data cleaning with a suitable example.</div> <div>OR</div> <div>ii. Given the following data (in increasing order) for the attribute age: 13, 15, 16, 16, 19, 20, 20, 21, 22, 22, 25, 25, 25, 25, 30, 33, 33, 35, 35, 35, 35, 36, 40, 45, 46, 52, 70.</div> <div><div>a. Use min-max normalization to transform the value 35 for age onto the range [0.0,1.0].</div><div>b. Use z-score normalization to transform the value 35 for age, where the standard deviation of age is 12.94 years.</div><div>c. Use normalization by decimal scaling to transform the value 35 for age.</div></div>	<div>[05]</div> <div>[05]</div>												



Q3 (b)	<p>i. Explain the relation between arithmetic, geometric and harmonic mean with suitable example.</p> <p style="text-align: center;">OR</p> <p>ii. Explain the various absolute measure of dispersion.</p> <p>iii. Find the range and coefficient of range of the following: 25, 67, 48, 53, 18, 39, 44.</p>	<p>[10]</p> <p>[05]</p> <p>[05]</p>																																
Q4 (a)	<p>1. Answer the following: Calculate the Pearson correlation coefficient for the Marks obtained by 5 students in algebra and trigonometry as given below:</p> <table border="1"><tr><td>Science</td><td>16</td><td>15</td><td>12</td><td>10</td><td>8</td></tr><tr><td>Geometry</td><td>11</td><td>18</td><td>10</td><td>20</td><td>17</td></tr></table> <p>ii. Find the SNR for the following data: 6, 24, 6, 14, 10.</p> <p>iii. What would be the standard deviation of the data if the mean is 45 and SNR is 39.5?</p> <p style="text-align: center;">OR</p> <p>2. Compute the Spearman's rank correlation coefficient of the following data. Also specify if the ρ value is perfect association of rank or no.</p> <table border="1"><tr><td>Rank in History (X)</td><td>3</td><td>5</td><td>1</td><td>6</td><td>7</td><td>2</td><td>8</td><td>9</td><td>4</td></tr><tr><td>Rank in Geography (Y)</td><td>5</td><td>3</td><td>2</td><td>6</td><td>8</td><td>1</td><td>7</td><td>9</td><td>4</td></tr></table>	Science	16	15	12	10	8	Geometry	11	18	10	20	17	Rank in History (X)	3	5	1	6	7	2	8	9	4	Rank in Geography (Y)	5	3	2	6	8	1	7	9	4	<p>[03]</p> <p>[02]</p> <p>[02]</p> <p>[07]</p>
Science	16	15	12	10	8																													
Geometry	11	18	10	20	17																													
Rank in History (X)	3	5	1	6	7	2	8	9	4																									
Rank in Geography (Y)	5	3	2	6	8	1	7	9	4																									
Q4.(b)	<p>i. Given the following data (in increasing order) for the attribute age: 4, 8, 9, 15, 21, 21, 24, 25, 26, 28, 29, 34. Calculate the following using a bin size of 3:</p> <p>a) Use smoothing by bin mean to smooth these data.</p> <p>b) Use smoothing by bin boundary to smooth these data.</p>	<p>[08]</p>																																
Q5 (a)	<p>Solve any two.</p> <p>i. Explain with an example the difference between correlation and causation.</p> <p>ii. Explain various Multivariate methods in feature engineering.</p> <p>iii.If two variables have a high correlation with a third variable, does this convey they will also be highly correlated? Is it even possible that A and B are positively correlated to another variable C? Is it possible that A and B are negatively correlated with each other?</p> <p>iv.Explain extreme values and limits in examining the relationship in correlation.</p>	<p>[05]</p> <p>[05]</p> <p>[05]</p> <p>[05]</p>																																
Q5 (b)	<p>i. Calculate the correlation coefficient for the following height (in inches) of mother (X) and their daughter (Y) using arbitrary origin method of Karl Pearson's coefficient.</p> <table border="1"><tr><td>X</td><td>65</td><td>66</td><td>67</td><td>67</td><td>68</td><td>69</td><td>70</td><td>72</td></tr><tr><td>Y</td><td>67</td><td>68</td><td>65</td><td>68</td><td>72</td><td>72</td><td>69</td><td>71</td></tr></table>	X	65	66	67	67	68	69	70	72	Y	67	68	65	68	72	72	69	71	<p>[05]</p>														
X	65	66	67	67	68	69	70	72																										
Y	67	68	65	68	72	72	69	71																										